



# Abstraction Pathologies In Markov Decision Processes

Manel Tagorti, Bruno Scherrer, Olivier Buffet, Joerg Hoffmann

## ► To cite this version:

Manel Tagorti, Bruno Scherrer, Olivier Buffet, Joerg Hoffmann. Abstraction Pathologies In Markov Decision Processes. ICAPS'13 workshop on Heuristics and Search for Domain-independent Planning (HSDIP), Jun 2013, Rome, Italy. hal-00907315

**HAL Id: hal-00907315**

**<https://inria.hal.science/hal-00907315>**

Submitted on 21 Nov 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Abstraction Pathologies In Markov Decision Processes

Manel Tagorti and Bruno Scherrer and Olivier Buffet

INRIA

Nancy, France

{manel.tagorti,bruno.scherrer,olivier.buffet}@inria.fr

Jörg Hoffmann

Saarland University

Saarbrücken, Germany

hoffmann@cs.uni-saarland.de

## Abstract

Abstraction is a common method to compute lower bounds in classical planning, imposing an equivalence relation on the state space and deriving the lower bound from the quotient system. It is a trivial and well-known fact that refined abstractions can only improve the lower bound. Thus, when we embarked on applying the same technique in the probabilistic setting, our firm belief was to find the same behavior there. We were wrong. Indeed, there are cases where *every* direct refinement step (splitting one equivalence class into two) yields strictly *worse* bounds. We give a comprehensive account of the issues involved, for two wide-spread methods to define and use abstract MDPs.

## Introduction

In classical planning, an abstraction is a mapping  $\alpha$  from the set of all states into a smaller set of abstract states ( $[s]_\alpha$  denoting the set of states  $t$  where  $\alpha(t) = \alpha(s)$ ). This is used to derive a lower bound  $h^\alpha(s)$  on the remaining cost of any state  $s$ . Namely,  $\alpha$  induces an abstract planning problem over the abstract state space: (i) an abstract state  $[s]_\alpha$  is a goal iff it contains at least one original goal state, and (ii) a transition from  $[s]_\alpha$  to  $[s']_\alpha$  exists iff there exist  $t \in [s]_\alpha$  and  $t' \in [s']_\alpha$  so that the original state space has a transition from  $t$  to  $t'$ . The abstract planning problem is a relaxed version of the original one –or, conversely, the original problem is more constrained– so that, given a state  $s$ , the cost of an abstract plan starting from  $[s]_\alpha$  is at most equal to the cost of a plan starting from  $s$ , and can thus be used as a lower bound  $h^\alpha(s)$ . Prominent examples of this method are pattern databases (Edelkamp 2001; Haslum et al. 2007) and merge-and-shrink abstractions (Helmert, Haslum, and Hoffmann 2007; Nissim, Hoffmann, and Helmert 2011; Katz, Hoffmann, and Helmert 2012).

A refinement of  $\alpha$  is an abstraction  $\alpha'$  resulting from  $\alpha$  by splitting some of the block states, i.e., for all  $s$  we have  $[s]_{\alpha'} \subseteq [s]_\alpha$ . It is commonplace that refinements can only improve the heuristic, i.e.,  $h^\alpha(s) \leq h^{\alpha'}(s)$  for all  $s$ : If we split block states apart, then the solution paths can only get longer and thus more costly (assuming non-negative costs as usual). Indeed, this observation is so simple that, to our knowledge, no-one yet bothered to state it in a paper and its first appearance is in Malte Helmert’s 2010 lecture slides.<sup>1</sup>

Our initial agenda in this research was to solve MDPs using heuristic search methods like LRTDP (Bonet and Geffner 2003), our focus being to compute heuristic functions by starting with a coarse abstraction and iteratively refining it. Against the background described above, as a warm-up exercise we embarked on proving that the essential property of refinements – they can only improve the heuristic – is true in that setting as well. Which was all fine, except we ended up proving the opposite.

First things first, to conduct this kind of research for MDPs one needs to first define what the “quotient system” and the corresponding heuristic functions are. This is non-trivial because, in difference to the classical case where all we are interested in is which states can transition to which other states in principle, now we need to define transition *probabilities* for the abstract MDP. To illustrate, if action  $a$  maps  $s$  into a state from  $[s']_\alpha$  with probability 0.9, but maps  $t \in [s]_\alpha$  into a state from  $[s']_\alpha$  with probability 0.1, which probability should we assign for  $a$  to map  $[s]_\alpha$  into  $[s']_\alpha$ ?

A simple answer is to assign the average probability over all states in  $[s]_\alpha$ . In the example, this would yield the transition probability 0.5. A main issue with this approach is that the resulting heuristic function – the value function of the abstract MDP – is neither a lower bound nor an upper bound on the value function of the original MDP. Givan et al. (2000) fix this by basically considering intervals of transition probabilities (in the example, the interval  $[0.1, 0.9]$ ). They derive a lower bound on the original value function (expected reward) by selecting the probabilities pessimistically, and derive an upper bound of the original value function by selecting the probabilities optimistically.

We first proved that, for the average-probability approach, there exist an MDP, a state  $s$ , an abstraction  $\alpha$ , and refined  $\alpha'$  so that the error of  $h^{\alpha'}(s)$  relative to the original value function is larger than that of  $h^\alpha(s)$ . This may be duly understood as an accident pertaining to the sketchy nature of this approach; indeed, as we show, the original MDP does not have to be non-deterministic to provoke this kind of behavior. However, we next proceeded to prove the same property for Givan et al.’s approach. Worse, even: We constructed an MDP, a state  $s$ , and  $\alpha$  so that *all* direct refinements  $\alpha'$  (resulting from  $\alpha$  by splitting a single block state) result in

<sup>1</sup><http://www.informatik.uni-freiburg.de/>

[~ki/teaching/ss10/aip/aip10.pdf](#)

strictly worse bounds. More naturally than for the average-probability approach, non-determinism is required for this, i.e., if the original MDP is deterministic then every refinement results in better bounds.

In the remainder of the paper, we first give the necessary background definitions. We then explain our results in the order described above.

## Markov Decision Processes and Abstractions

Markov Decision Processes (MDPs) are a general framework for modeling decision-making problems in stochastic environments. We define an MDP as follows

**Definition 1.** A Markov Decision Process is given by a (finite) state space  $S$ , a (finite) action space  $A$ , a reward function  $R : S \times A \rightarrow \mathbb{R}$  and transition probabilities  $p(s, a, s')$  which determine the probabilities of transition when performing action  $a$  in state  $s$ .

A deterministic policy  $\pi : S \rightarrow A$  assigns an action to each state, and we are looking for the optimal policy  $\pi^*$ , i.e., one that maximizes for all  $s \in S$  the return

$$V^\pi(s) = E_\pi \left[ \sum_{t=0}^{\infty} \gamma^t R_t | s_0 = s \right],$$

where  $\gamma$  is the discount rate taken in  $[0, 1)$ .

The related value  $V^*$  is the unique solution of the equation  $V^* = TV^*$  where  $T$  is the Bellman operator defined as :

$$\begin{aligned} \forall s, V^*(s) &= TV^*(s) \\ &= \max_a R(s, a) + \gamma \sum_{s' \in S} p(s, a, s') V^*(s'). \end{aligned}$$

Determining  $\pi^*$  may be infeasible in MDPs with large state spaces. In this paper we simplify the problem by employing state abstractions. Abstractions provide a smaller representation  $M_\alpha$  of the original MDP  $M$ . The image of  $M$  under an abstraction  $\alpha$  is defined on a state space  $S_\alpha$  smaller than  $S$ . Indeed  $S_\alpha$  is a partition of  $S$  consisting of block states  $[s]_\alpha$ . We assign to each block-state, given an action  $a \in A$ , the reward  $R([s]_\alpha, a)$  and the transition probabilities  $p([s]_\alpha, a, [s_1]_\alpha)$  for all  $[s_1]_\alpha \in S_\alpha$ . The useful abstractions are the ones that induce a small error of approximation when considering  $M_\alpha$  instead of  $M$ . We would like to identify such abstractions by comparing a given abstraction to its (direct) refinement, in terms of approximation error.

## Abstractions' Refinement

**Definition 2.** Let  $\alpha$  and  $\alpha'$  be two abstractions of an MDP  $M$ . We say that  $\alpha'$  is finer than  $\alpha$ , denoted  $\alpha' \succeq \alpha$ , iff for any states  $s, s' \in S$ ,  $\alpha'(s) = \alpha'(s')$  implies  $\alpha(s) = \alpha(s')$ . We can also say that  $\alpha$  is coarser than  $\alpha'$ , denoted  $\alpha \preceq \alpha'$ .

We have  $\alpha'$  a direct refinement of  $\alpha$  if there exist states  $s_1, s_2 \in S$  such that  $[s_1]_\alpha = [s_2]_{\alpha'}$ ,  $[s_1]_{\alpha'} \neq [s_2]_{\alpha'}$ ,  $[s_1]_\alpha = [s_1]_{\alpha'} \cup [s_2]_{\alpha'}$ , and  $\alpha'(s) = \alpha(s)$  for all  $s \in S \setminus [s_1]_\alpha$ .

We show in what follows that the error induced by  $\alpha'$  may in some cases be higher than the one induced by  $\alpha$ . But before that we have first to specify the parameters (rewards and transition probabilities) related to the abstract representation  $M_\alpha$  of  $M$ .

## Average MDPs

We consider in this section the abstraction  $\alpha$  that connects an MDP  $M$  to its average representation. In other words,  $\alpha$  maps an MDP  $M$  defined on  $S$  to an average MDP  $M_\alpha$ , defined on  $S_\alpha$ , and admits as parameters, the averages of rewards and transitions over all states contained in the block state  $[s]_\alpha$  ((Ortner 2011)), i.e., we have for all  $a \in A$ :

$$\begin{aligned} R([s]_\alpha, a) &= \frac{1}{|[s]_\alpha|} \sum_{s_1 \in [s]_\alpha} R(s_1, a) \text{ and} \\ p([s]_\alpha, a, [s']_\alpha) &= \frac{1}{|[s]_\alpha|} \sum_{s_1 \in [s]_\alpha} \sum_{s_2 \in [s']_\alpha} p(s_1, a, s_2). \end{aligned}$$

We denote here by  $|[s]_\alpha|$  the cardinal of all states in  $[s]_\alpha$ .

We choose as approximation error  $E_\alpha$ , the average error, estimated by taking the average difference between the true value  $V$  taken in a state  $s$  and the value  $V_\alpha$  of its corresponding block state  $[s]_\alpha$ ,

$$E_\alpha = \frac{1}{|S|} \sum_{s \in S} |V_\alpha([s]_\alpha) - V(s)|.$$

This approximation error may increase when we refine the abstraction  $\alpha$ . We illustrate in Figure 1 an example in which the number of states where the (local) error increases-after a refinement-is greater than the one where the (local) error decreases, causing the increase of the average error.

**Proposition 1.** There exists a deterministic MDP  $M$ , an abstraction  $\alpha$  and a refinement  $\alpha'$  of  $\alpha$  such that  $E_\alpha < E_{\alpha'}$ .

*Proof.* Consider the MDP  $M$  in Figure 1 with a single action  $\{a\}$  and a discount rate  $\gamma = 1$  (the result would not change for  $\gamma < 1$  but close enough to 1). The states in  $\{2, \dots, k\}$  ( $k > 2$ ) are similar: they admit the same rewards ( $R = 0$ ) and have the same dynamics : they reach the neighboring state with probability one. The states in  $\{k+1, \dots, n\}$  are also similar: they all reach the goal  $G$  with probability one and they admit a non-negative reward  $R_2$ . The state 1 admits a reward  $R(1) = R_1 \ll R_2$  and reaches the goal with probability one.

Taking  $V(G) = 0$ , we then have  $V(1) = R_1$ , and  $V(i) = R_2$  for all  $i$  in  $\{2, \dots, n\}$ .

Based on those similarities one can construct a perfectly suitable abstraction  $\alpha_0$  ( $E_{\alpha_0} = 0$ ) which aggregates similar states in the same block, i.e., in our case  $\alpha_0 : S \rightarrow 1, \{2, \dots, k\}, \{k+1, \dots, n\}, G$ .

Consider now the abstraction  $\alpha_1 : S \rightarrow \{1, \dots, k\}, \{k+1, \dots, n\}, G$ , where the states 1 and  $\{2, \dots, k\}$  are in the same block (Figure 1). The similarity will be then broken resulting in a strictly positive error  $E_{\alpha_1}$ . The related values are

$$\begin{aligned} V_{\alpha_1}(\{k+1, \dots, n\}) &= \frac{(n-k)R_2}{n-k} = R_2 \text{ and} \\ V_{\alpha_1}(\{1, \dots, k\}) &= \frac{R_1}{k} + \frac{k-2}{k} V_{\alpha_1}(\{1, \dots, k\}) + \\ &\quad \frac{1}{k} V_{\alpha_1}(\{k+1, \dots, n\}) \\ &= \frac{R_1 + R_2}{2} \neq R_2. \end{aligned}$$

And the induced error is

$$E_{\alpha_1} = \frac{1}{n} \sum_{i=1}^k |V(i) - V_{\alpha_1}(\{1, \dots, k\})| \sim k \frac{R_2}{2n} \text{ for } R_1 \ll R_2.$$

Let  $\alpha_2 : S \rightarrow \{1, \dots, n\}, G$  be the abstraction which aggregates states  $1, 2, \dots, n$  together (Figure 1), the value  $V_{\alpha_2}(\{1, \dots, n\})$  is equal to

$$V_{\alpha_2}(\{1, \dots, n\}) = \frac{R_1 + (n-k)R_2}{n} + \frac{k-1}{n} V_{\alpha_2}(\{1, \dots, n\}) \\ \sim R_2 \text{ for } k \ll n$$

and hence

$$E_{\alpha_2} \sim \frac{1}{n} |V(1) - V_{\alpha_2}(\{1, \dots, n\})| \sim \frac{1}{n} R_2.$$

We can see that the error  $E_{\alpha_1}$  is strictly larger than  $E_{\alpha_2}$  for a number of states  $k$  strictly greater than 2.  $\square$

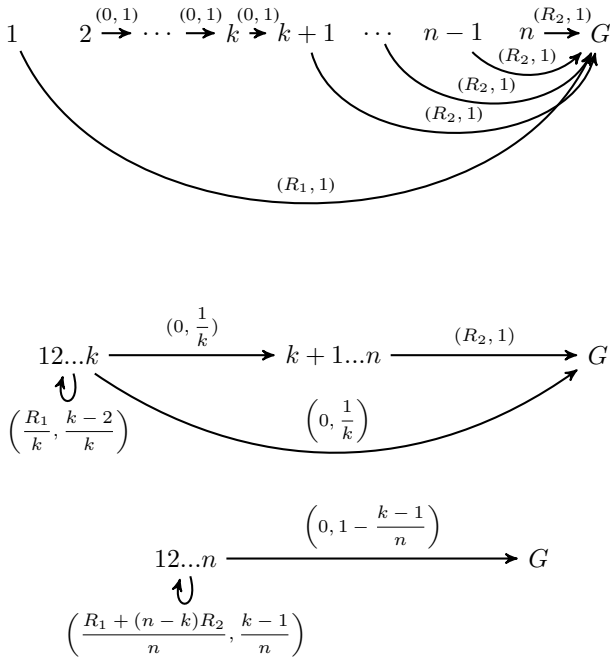


Figure 1: From the top to the bottom: The MDP  $M$ , the abstraction  $\alpha_1$  and the abstraction  $\alpha_2$  ( $\alpha_0$  is not shown). The parentheses denote the reward (on the left) and the transition probabilities (on the right).

### Bounded-parameter MDPs

Rather than approaching  $V^*$  with fixed values  $V_\alpha$ , it is possible to establish bounds on the MDP's parameters so as to get bounds on  $V^*$ . So we considered the abstraction  $\alpha$  that associates, for each (action)  $a$  in  $A$ , each block state  $[s]_\alpha$  with

the intervals' parameters ((Givan, Leach, and Dean 2000)):

$$R^\dagger([s]_\alpha, a) = [\min_{s \in [s]_\alpha} R(s, a), \max_{s \in [s]_\alpha} R(s, a)] \\ p^\dagger([s]_\alpha, a, [s']_\alpha) = [\min_{s_1 \in [s]_\alpha} \sum_{s_2 \in [s']_\alpha} p(s_1, a, s_2), \\ \max_{s_1 \in [s]_\alpha} \sum_{s_2 \in [s']_\alpha} p(s_1, a, s_2)].$$

We get hence what we call a *bounded parameter Markov Decision process* (BMDP) more commonly defined as:

**Definition 3.** A Bounded parameter Markov Decision Process is given by a (finite) state space  $\Sigma$ , a (finite) action space  $A$ , an interval of rewards  $R^\dagger(\sigma, a)$ ,  $\forall \sigma \in \Sigma, a \in A$ , and an interval of transition probabilities  $p^\dagger(\sigma, a, \sigma')$ ,  $\forall \sigma, \sigma' \in \Sigma, a \in A$ .

Each state of the BMDP has a range of values depending on  $R$  and  $p$ . We can assign to each state a closed interval of value functions  $[V^-(\sigma), V^+(\sigma)]$  where  $V^-$  corresponds to the pessimistic bound and  $V^+$  to the optimistic one.

Our abstract representation  $M_\alpha$  is then a BMDP where the state space  $\Sigma$  coincides with state space  $S_\alpha$ . We would like to estimate the value bounds  $V_\alpha^+$  and  $V_\alpha^-$  related to  $M_\alpha$ . Givan *et al.* have proposed an algorithm, the interval value iteration IVI, to do so. To make explicit the two Bellman's operators hidden behind this algorithm ( $T^+$  and  $T^-$ ), we first need to introduce the notion of *compatibility* with respect to an abstraction.

**Definition 4.** An MDP  $N$  is compatible with  $M$  with respect to the abstraction  $\alpha$  if for all  $s$ , and for all  $a$ ,  $R_N(s, a) \in R_M^\dagger([s]_\alpha, a)$  and for all  $s, s', a$ ,  $\sum_{s_1 \in [s']_\alpha} p_N(s, a, s_1) \in p_M^\dagger([s]_\alpha, a, [s']_\alpha)$ . The set of all MDPs compatible with  $M$  with respect to  $\alpha$  is denoted  $[M]_\alpha$ .

Figure 2 gives an example of an MDP  $N$  compatible with an MDP  $M$  with respect to the abstraction  $\alpha : 1, 2, 3 \rightarrow \{1, 2\}, 3$ .

We claim that the Bellman operators  $T^+$  and  $T^-$  used to estimate  $V_\alpha^+$  and  $V_\alpha^-$  can be written in this way, for all  $s$ :

$$T^+[V](s) = \max_{a \in A} \max_{N \in [M]_\alpha} R_N(s, a) + \gamma \sum_{s' \in S} p_N(s, a, s') V(s'),$$

$$T^-[V](s) = \max_{a \in A} \min_{N \in [M]_\alpha} R_N(s, a) + \gamma \sum_{s' \in S} p_N(s, a, s') V(s').$$

By taking iteratively the max (resp the min) on the set of compatible MDPs and by choosing the optimal policy, we can see that those two operators converge to fixed values. Indeed  $T^+$  and  $T^-$  are  $\gamma$ -contracting so, by the Banach fixed point Theorem, they admit unique fixed points  $V_\alpha^+$  and  $V_\alpha^-$  (which are constant per block), i.e.,  $T^+ V_\alpha^+ = V_\alpha^+$  and  $T^- V_\alpha^- = V_\alpha^-$ . There exists an optimistic MDP  $M^{opt}$  (respectively a pessimistic MDP  $M^{pes}$ ) and a corresponding optimal (optimistic) policy  $\pi^{opt}$  (respectively an optimal (pessimistic) policy  $\pi^{pes}$ ) for which the value  $V_\alpha^+$  (resp  $V_\alpha^-$ ) is reached. The MDPs  $M^{opt}$  and  $M^{pes}$  belong to  $[M]_\alpha$ .

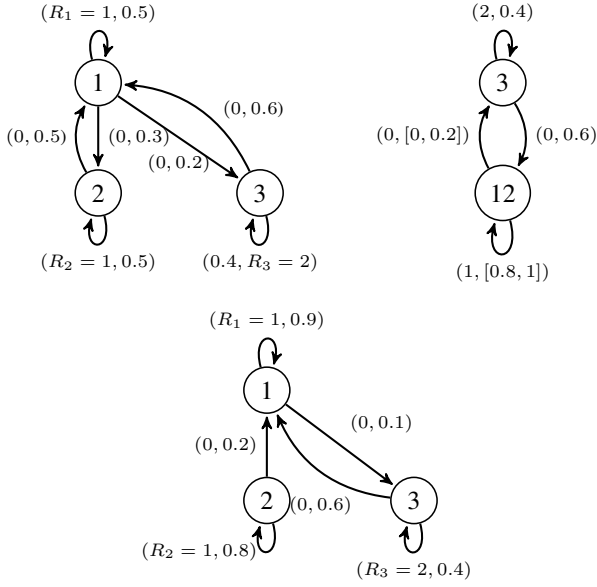


Figure 2: From left to right, the MDP  $M$ , the BMDP  $M_\alpha$ ,  $\alpha : 1, 2, 3 \rightarrow \{1, 2\}, 3$  and the MDP  $N$  compatible with  $M$  with respect to  $\alpha$ . The parenthesis denote the reward (on the left) and the transitions (on the right).

Before proceeding to the next step –the choice of the abstraction that would ensure better bounds–, we first show that these values are indeed bounds on  $V^*$ . This is precisely what is stated in this following theorem.

**Theorem 1.** (Givan, Leach, and Dean 1997) *For any MDP  $M$  and abstraction  $\alpha$  of the states of  $M$ , bounds on the BMDP  $M_\alpha$  apply also to  $M$ , i.e.,  $\forall s \in S$ ,  $V^*(s) \in [V_\alpha^-(s)_\alpha, V_\alpha^+(s)_\alpha]$ .*

*Proof.*<sup>2</sup> The proof is done using Value Iteration. With initial value  $V^0 = V_\alpha^+$  we have, for all  $s$ ,

$$V^1(s) = \max_a R_M(s, a) + \gamma \sum_{s' \in S} p_M(s, a, s') V_\alpha^+([s']_\alpha).$$

Since  $V_\alpha^+$  is a fixed point of  $T^+$  and  $M \in [M]_\alpha$ , we can see that

$$\begin{aligned} V^1(s) &\leq \max_{a \in A} \max_{N \in [M]_\alpha} R_N(s, a) + \gamma \sum_{s' \in S} p_N(s, a, s') V_\alpha^+([s']_\alpha) \\ &\leq V_\alpha^+([s]_\alpha) = V^0(s). \end{aligned}$$

The Bellman operator  $T$  is monotone, so that we have  $T^n V^1 \leq T^n V^0 = V^0$ . By taking the limit, we get  $V^*(s) \leq V_\alpha^+([s]_\alpha)$ . A similar proof may be applied to the pessimistic bound.  $\square$

<sup>2</sup>This proof does not appear in (Givan, Leach, and Dean 1997) but it has been established thanks to a correspondence with R. Givan.

## Value Bounds using finer abstractions

As previously done for the average model, we will compare the two errors  $G_\alpha$  and  $G_{\alpha'}$  induced by the BMDP  $M_\alpha$  and its direct refinement  $M_{\alpha'}$ , where, for an abstraction  $\alpha$ ,  $G_\alpha$  measures the gap between the two bounds in each state, for all  $s$ :

$$G_\alpha(s) = (V_\alpha^+([s]_\alpha) - V_\alpha^-([s]_\alpha)).$$

The proposition below states a sufficient condition under which the finer abstraction  $\alpha'$  yields better value bounds than  $\alpha$  and therefore decreases the error  $G_\alpha$ . Indeed, if the set of MDPs compatible with  $M$  with respect to the finer abstraction are also compatible with  $M$  with respect to the coarser abstraction, then we get the inclusion of the value function intervals.

**Proposition 2.** *Given an MDP  $M$  and two abstractions  $\alpha$  and  $\alpha'$  s.t.  $\alpha' \succeq \alpha$  and  $[M]_{\alpha'} \subseteq [M]_\alpha$  then, for all  $s \in S$ ,*

$$[V_{\alpha'}^-([s]_{\alpha'}), V_{\alpha'}^+([s]_{\alpha'})] \subseteq [V_\alpha^-([s]_\alpha), V_\alpha^+([s]_\alpha)].$$

*Proof.* Value Bounds computation can be considered as a case of an alternating two players stochastic game where one choose the optimal policy while the other choose the optimal MDP. For the upper bound, we can then invert the two max in the Bellman operators' expressions ((Givan, Leach, and Dean 2000)). This would not change the final result, so we have

$$V_\alpha^+([s]_\alpha) = \max_{N \in [M]_\alpha} V_N^*(s).$$

For the lower bound, more detailed arguments have been established in (Bertsekas and Tsitsiklis 1996) to set the inversion of the max and the min terms, so we get:

$$V_\alpha^-([s]_\alpha) = \min_{N \in [M]_\alpha} V_N^*(s).$$

Given that  $[M]_{\alpha'}$  is included in  $[M]_\alpha$ , then by taking the max, (respectively the min) the result follows.  $\square$

We will study next two cases : the deterministic case, where the sufficient condition is always satisfied, and the stochastic case, where it is not necessarily satisfied, for which we will give an example.

### Deterministic case: probabilities in $\{0, 1\}$

**Corollary 1.** *If we consider a deterministic MDP and two abstractions  $\alpha$  and  $\alpha'$ , where  $\alpha'$  is a direct refinement of  $\alpha$ , then for all  $s \in S$ ,*

$$[V_{\alpha'}^-([s]_{\alpha'}), V_{\alpha'}^+([s]_{\alpha'})] \subseteq [V_\alpha^-([s]_\alpha), V_\alpha^+([s]_\alpha)].$$

*Proof.* Let us consider an MDP  $N$  compatible with  $M$  under  $\alpha'$ . We will show that the sufficient condition of Proposition 2 is satisfied:  $N$  is also compatible with  $M$  under  $\alpha$ . Note that having  $N$  compatible with  $M$  with respect to an abstraction  $\alpha$ , according to Definition 4, is equivalent to having the inclusion of the parameter intervals i.e., for all  $s$ , and for all  $a$ ,  $R_N^\dagger([s]_\alpha, a) \subseteq R_M^\dagger([s]_\alpha, a)$  and for all  $s, s', a$ ,  $p_N^\dagger([s]_\alpha, a, [s']_\alpha) \subseteq p_M^\dagger([s]_\alpha, a, [s']_\alpha)$ . Let  $[s_1]_\alpha$  be the state block in  $S_\alpha$  that we split into two blocks  $[s_2]_{\alpha'}$  and  $[s_3]_{\alpha'}$  (i.e., we have  $[s_2]_{\alpha'} \cup [s_3]_{\alpha'} = [s_1]_\alpha$ ). It is easy to

check the inclusion of reward intervals under  $\alpha$  as for each action  $a$ :

$$\min_{s \in [s_1]_\alpha} R_N(s, a) = \min(\min_{s \in [s_2]_{\alpha'}} R_N(s, a), \min_{s \in [s_3]_{\alpha'}} R_N(s, a)).$$

Also, using the  $\alpha'$  compatibility hypothesis we have

$$\min_{s \in [s_i]_{\alpha'}} R_N(s, a) \geq \min_{s \in [s_i]_{\alpha'}} R_M(s, a) \text{ for } i \text{ in } \{2, 3\}$$

then

$$\min_{s \in [s_1]_\alpha} R_N(s, a) \geq \min_{s \in [s_1]_\alpha} R_M(s, a).$$

By reasoning in a similar way for the upper bound, we get the inclusion of the reward intervals. The same arguments may be employed to state the inclusion of outgoing transition probabilities  $p([s_1]_\alpha, a, [s_4]_\alpha)$  for all  $s_4$  in  $S$  and  $a$  in  $A$ . So we will mainly focus on ingoing transition probabilities  $p([s_5]_\alpha, a, [s_1]_\alpha)$  for all  $s_5$  in  $S$  and we will show that

$$\min p_N([s_5]_\alpha, a, [s_1]_\alpha) \geq \min p_M([s_5]_\alpha, a, [s_1]_\alpha). \quad (1)$$

Since we work in a deterministic environment, probabilities can only take the values 0 or 1. For the case where  $\min p_N([s_5]_\alpha, a, [s_1]_\alpha) = 1$ , inequality (1) is always verified. Now,  $\min p_N([s_5]_\alpha, a, [s_1]_\alpha) = 0$  implies that there exist a state  $s'$  in  $[s_5]_\alpha$  and a block state  $[s_6]_\alpha$  distinct from  $[s_1]_\alpha$  such that  $p_N(s', a, [s_6]_\alpha) = 1$ . So we can find a state  $s''$  in  $[s_5]_\alpha$  such that  $p_M(s'', a, [s_6]_\alpha) = 1$  as  $N$  is compatible with  $M$  under  $\alpha'$ . We then have Equation (1) and by Proposition 2 the final result follows.  $\square$

**Stochastic case** We would like to have an equivalent of Proposition 2 for the stochastic case but it turns out that in general the sufficient condition is no more fulfilled. In fact when we move to stochastic transitions, the sufficient condition in Proposition 2 becomes harder to satisfy. The successful MDP  $M$  has to verify specific conditions that depend on the choice of the abstractions  $\alpha$  and  $\alpha'$ . In other words, given an MDP  $M$  and an abstraction  $\alpha$  there does not always exist a refined abstraction  $\alpha'$  such that the sufficient condition is satisfied. Figure 3 shows a model of MDP in which we can not find the appropriate direct refinement  $\alpha'$ . Three states which admit the same rewards  $R(1) = R(2) = R(3) = 1$  and behave identically in the block  $\{1, 2, 3\}$  (the same probabilities in regards to the block  $\{1, 2, 3\}$ ,  $p(1, \{1, 2, 3\}) = p(2, \{1, 2, 3\}) = p(3, \{1, 2, 3\}) = 0.7$ ). States 4 and 5 are goal states ( $V(4) = V(5) = 0$ ). We can find an MDP  $N$  compatible with  $M$  with respect to the abstraction  $\alpha'$  ( $S \rightarrow \{1, 2\}, 3, 4, 5$ ) but not compatible with  $M$  with respect to the abstraction  $\alpha$  ( $S \rightarrow \{1, 2\}, 3, 4, 5$ ). It suffices to take an MDP  $N$  whose parameter intervals are included in  $M$ 's parameter intervals under  $\alpha'$  but admits real intervals rather than fixed real values under  $\alpha$ .

The condition stated in Proposition 2 is a sufficient condition, we can not *a priori* conclude about the existence of a refinement that would make the error of approximation decrease. Nevertheless, we have identified models of MDPs where every direct refinement strictly increases the error.

**Proposition 3.** *There exists an MDP  $M$  and an abstraction  $\alpha$  such that, for any direct refinement  $\alpha'$  of  $\alpha$ , we have: (1) for all  $s$  in  $S$ ,  $G_{\alpha'}(s) \geq G_\alpha(s)$ , and (2) there exists  $s$  in  $S$  where  $G_{\alpha'}(s) > G_\alpha(s)$ .*

*Proof.* The MDP shown in Figure 3 is an example of such a model. The value  $V_\alpha(\{1, 2, 3\})$  related to the block  $\{1, 2, 3\}$  is a perfect heuristic for the states 1, 2 and 3. In fact we can see (for  $\gamma = 1$ ), using Theorem 1, that the optimistic and pessimistic bounds coincide:  $V_\alpha^-(\{1, 2, 3\}) = V_\alpha^+(\{1, 2, 3\}) = V(1) = V(2) = V(3) = 3.33$ . If we refine the block  $\{1, 2, 3\}$  by splitting the block  $\{1, 2, 3\}$  into  $\{1, 2\}$  and  $\{3\}$ , then the error of approximation will strictly increase and we will obtain an interval of values rather than a precise value given by the coarser abstraction (we got those values by using Givan et al.'s algorithm). The approximation error  $G_\alpha$  equals to 0 for all the states while the approximation error  $G_{\alpha'}$  is:

$$G_{\alpha'}(1) = G_{\alpha'}(2) \simeq 1 > 0 = G_\alpha(1) = G_\alpha(2) \text{ and } G_{\alpha'}(3) \simeq 1 > G_\alpha(3) = 0.$$

Only one of all the possible refinements is detailed here but the same happens for the two other direct refinements:

- For the abstraction

$$\alpha' : 1, 2, 3, 4, 5 \rightarrow \{1, 3\}, \{2\}, \{4\}, \{5\}$$

$$\text{we have } V_{\alpha'}^\pm(\{1, 3\}) = [2.5, 4.78] \text{ and } V_{\alpha'}^\pm(\{2\}) = [2.75, 4.34].$$

- For the abstraction

$$\alpha' : 1, 2, 3, 4, 5 \rightarrow \{1\}, \{2, 3\}, \{4\}, \{5\}$$

$$\text{we have } V_{\alpha'}^\pm(\{2, 3\}) = [2.25, 5.29] \text{ and } V_{\alpha'}^\pm(\{1\}) = [2.9, 4.11].$$

$\square$

We can even have a model inspired from the one above where the error strictly increases in each state. Indeed by taking  $\gamma = 0.9$  and by changing the transition probabilities in the initial model  $M$  to  $p(4, 4) = p(5, 5) = 0.9$  and  $p(4, 1) = p(5, 1) = 0.1$  (we keep the same rewards  $R(4) = R(5) = 0$ ), we can see that the gap  $G_{\alpha'}$  is strictly higher than  $G_\alpha = 0$  for each direct refinement  $\alpha'$ .

## Related Work

Our work shows the limitations of some abstractions (state abstractions) in which refining an abstraction may increase the approximation error. This has already been observed in (Waugh et al. 2009) in the case of an  $i$ -player poker game ( $i$  greater than 2). They looked at the exploitability<sup>3</sup> of each player's strategy with respect to each abstraction and they showed that it may increase while considering finer card abstractions. The same phenomenon has been observed

<sup>3</sup>The exploitability is a metric connected to the Nash equilibrium strategy. It is equal to 0 in the case of a two-player extensive game.

when they considered betting abstractions – by restricting the number of betting options in each sequence of the game. Lately another example about ”action abstractions” pathologies has been provided in (Sandholm and Singh 2012).

It is important to notice here that, contrary to what has been done in those works, this paper deals with single action and player models. This suggests that, in some abstractions, notably BMDP abstraction, stochasticity alone can explain this pathological behavior, and that we do not need to consider the more general case of two-player game: the issues appear even in the case of a one-player game (MDP). Indeed, our counter-examples do not even contain any actual action choice, thus identifying this kind of pathology in a very canonical framework.

Interestingly, Kattenbelt et al. (2010) introduce a variant of abstraction for MDPs that, according to their results, does not exhibit said pathology: Every abstraction refinement step results in an improved bound. An interesting open question is how exactly their framework relates to BMDP abstraction.

## Conclusions

Somewhat surprisingly, refining an abstraction does not guarantee, in the MDP setting, a refined, i.e., better, approximation of the value function. From a practical perspective, this observation might be reasonably classified as ”odd but not crucial” – the loss of this guarantee is not per se an argument against trying to apply abstraction techniques for the computation of heuristic functions, as known from classical planning. The observation *might* be relevant to the practical effectiveness of such methods, where paying a higher price for the abstraction may result in less accuracy. But it remains to be seen whether that is of practical importance.

From a theoretical perspective, we believe that our observations could be of importance for a better understanding of the methods involved. In that regard, our investigation is but a small start into the subject matter. In particular, our vision was and is to identify sufficient criteria, in the bounded-parameter MDP setting, for the error to not increase. If such a criterion is efficiently testable, or can at least be reasonably well approximated, then it could serve as a well-informed guidance during the abstraction refinement process. For the moment, we don’t know how such a criterion could be designed. An interesting observation in this context is that our counter-example refines an abstraction that already is a bisimulation.<sup>4</sup> Does that tell us something about the general case? Another question is whether increasing the ”extent” of the per-step refinement helps (instead of splitting a single block-state, split 2, 3, ... block-states). Does there exist a non-trivial method (not refining all the way to the original MDP) that guarantees, for any MDP and abstraction thereof, the existence of a refinement step reducing the error? And

<sup>4</sup>If we aggregate the states 4 and 5 together, in the abstract representation  $\alpha$ , we get a bisimulation: all the states behave similarly in regards to each block of the partition  $\{1, 2, 3\}, \{4, 5\}$ . We have  $R(1) = R(2) = R(3) = 1$ ,  $p(1, \{1, 2, 3\}) = p(2, \{1, 2, 3\}) = p(3, \{1, 2, 3\}) = 0.7$  and  $p(1, \{4, 5\}) = p(2, \{4, 5\}) = p(3, \{4, 5\}) = 0.3$ , states 4 and 5 are goal states.

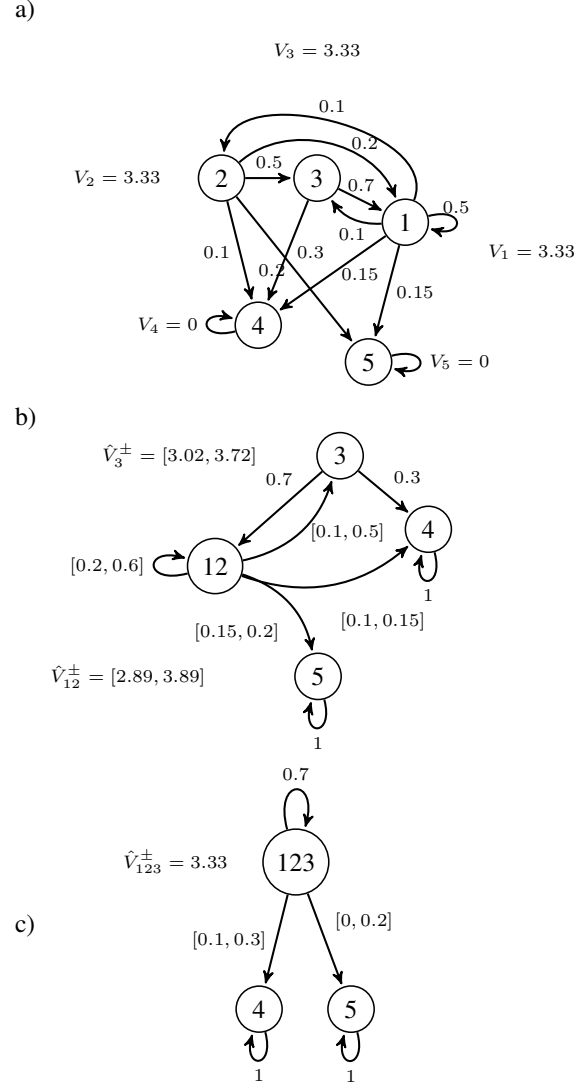


Figure 3: From top to bottom: the MDP  $M$ , the abstraction  $\alpha'$ , the abstraction  $\alpha$ . Edges are annotated with transition probabilities. The reward is:  $R(1) = R(2) = R(3) = 1 (= R(\{1, 2\}) = R(\{1, 2, 3\}))$ , and  $R(4) = R(5) = 0$ .

can that method be made practical? We believe these are interesting questions for future research, and hope other researchers might join us in exploring them.

**Acknowledgments.** We thank the anonymous HSDIP'13 reviewers, whose comments helped to improve the paper.

## References

- Bertsekas, D. P., and Tsitsiklis, J. 1996. *Neuro-Dynamic Programming*. Athena Scientific.
- Bonet, B., and Geffner, H. 2003. Labeled RTDP: Improving the convergence of real-time dynamic programming. In Giunchiglia, E.; Muscettola, N.; and Nau, D., eds., *Proceedings of the 13th International Conference on Automated Planning and Scheduling (ICAPS-03)*, 12–21. Trento, Italy: Morgan Kaufmann.
- Edelkamp, S. 2001. Planning with pattern databases. In Cesta, A., and Borrajo, D., eds., *Recent Advances in AI Planning, 6th European Conference on Planning (ECP-01)*, Lecture Notes in Artificial Intelligence, 13–24. Toledo, Spain: Springer-Verlag.
- Givan, R.; Leach, S. M.; and Dean, T. 1997. Model reduction techniques for computing approximately optimal solutions for Markov Decision Processes. *Artificial Intelligence* 122(1-2):1–8.
- Givan, R.; Leach, S. M.; and Dean, T. 2000. Bounded-parameter Markov Decision Processes. *Artificial Intelligence* 122(1-2):71–109.
- Haslum, P.; Botea, A.; Helmert, M.; Bonet, B.; and Koenig, S. 2007. Domain-independent construction of pattern database heuristics for cost-optimal planning. In Howe, A., and Holte, R. C., eds., *Proceedings of the 22nd National Conference of the American Association for Artificial Intelligence (AAAI-07)*, 1007–1012. Vancouver, BC, Canada: AAAI Press.
- Helmert, M.; Haslum, P.; and Hoffmann, J. 2007. Flexible abstraction heuristics for optimal sequential planning. In Boddy, M.; Fox, M.; and Thiebaux, S., eds., *Proceedings of the 17th International Conference on Automated Planning and Scheduling (ICAPS-07)*, 176–183. Providence, Rhode Island, USA: Morgan Kaufmann.
- Kattenbelt, M.; Kwiatkowska, M.; Norman, G.; and Parker, D. 2010. A game-based abstraction-refinement framework for Markov Decision Processes. *Formal Methods in System Design* 36(3):246–280.
- Katz, M.; Hoffmann, J.; and Helmert, M. 2012. How to relax a bisimulation? In Bonet, B.; McCluskey, L.; Silva, J. R.; and Williams, B., eds., *Proceedings of the 22nd International Conference on Automated Planning and Scheduling (ICAPS 2012)*. AAAI Press.
- Nissim, R.; Hoffmann, J.; and Helmert, M. 2011. Computing perfect heuristics in polynomial time: On bisimulation and merge-and-shrink abstraction in optimal planning. In Walsh, T., ed., *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI'11)*, 1983–1990. AAAI Press/IJCAI.
- Ortner, R. 2011. Adaptive aggregation for reinforcement learning in average reward Markov Decision Processes. *Annals of Operational Research*.
- Sandholm, T., and Singh, S. 2012. Lossy stochastic game abstraction with bounds. In *Proceedings of the 13th ACM Conference on Electronic Commerce, EC '12*, 880–897. ACM.
- Waugh, K.; Schnizlein, D.; Bowling, M.; and Szafron, D. 2009. Abstraction pathologies in extensive games. In *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems - Volume 2, AAMAS '09*, 781–788.